

Predicting heart disease using machine learning

This notebook looks into using various python-based machine learning and data science libraries in an attempt to build a machine learning model capable of predicting whether or not someone has heart disease based on their medical attributes.

we're going to take the following approach:

1. problem definition
2. data
3. Evaluation
4. features
5. experimentation

1. Problem Definition

In a statement,

Given clinical parameters about a patient, can we predict whether or not they have heart disease?

2. Data

originally from kaggle

3. Evaluation

If we can reach 95% accuracy at predicting whether or not a patient has heart disease during the proof of concept, we'll pursue the project

4. Features

This is where you'll get different information about each of the features in your data.

create data dictionary

1. age - age in years
2. sex - (1 = male; 0 = female)
3. cp - chest pain type *0: Typical angina: chest pain related decrease blood supply to the heart *1: Atypical angina: chest pain not related to heart *2: Non-anginal pain: typically esophageal spasms (non heart related) *3: Asymptomatic: chest pain not showing signs of disease
4. trestbps - resting blood pressure (in mm Hg on admission to the hospital) anything above 130-140 is typically cause for concern
5. chol - serum cholestoral in mg/dl serum = LDL + HDL + .2 * triglycerides above 200 is cause for concern fbs - (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false) '>126' mg/dL signals diabetes restecg - resting electrocardiographic results 0: Nothing to note 1: ST-T Wave abnormality can range from mild symptoms to severe problems signals non-normal heart beat 2: Possible or definite left ventricular hypertrophy Enlarged heart's main pumping chamber thalach - maximum heart rate achieved exang - exercise induced angina (1 = yes; 0 = no) oldpeak - ST depression induced by exercise relative to rest looks at stress of heart during exercise unhealthy heart will stress more slope - the slope of the peak exercise ST segment 0: Upsloping: better heart rate with exercise (uncommon) 1: Flatsloping: minimal change (typical healthy heart) 2: Downsloping: signs of unhealthy heart ca - number of major vessels (0-3) colored by flourosopy colored vessel means the doctor can see the blood passing through the more blood movement the better (no clots) thal - thalium stress result 1,3: normal 6: fixed defect: used to be defect but ok now 7: reversable defect: no proper blood movement when exercising target - have disease or not (1=yes, 0=no) (= the predicted attribute) Note: No personal identifiable information (PPI) can be found in the dataset.

It's a good idea to save these to a Python dictionary or in an external file, so we can look at them later without coming back here.

preparing the data

we're going to use pandas, matplotlib, and numpy for data analysis and manipulation

```
In [1]: 1 # import all the tools we need
2 # Regular EDA(exploratory data analysis) and plotting libraries
3
4 import pandas as pd
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # we want our plots to appear inside the notebook
10 %matplotlib inline
11
12 # models from scikit-Learn
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.ensemble import RandomForestClassifier
16
17 # model Evaluation
18 from sklearn.model_selection import train_test_split, cross_val_score
19 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
20 from sklearn.metrics import confusion_matrix, classification_report
21 from sklearn.metrics import precision_score, recall_score, f1_score
22 from sklearn.metrics import plot_roc_curve
```

```
In [2]: 1 import os
2 cwd = os.getcwd()
3 print(cwd)
```

C:\Users\USER\Desktop\heart-disease-project

Load data

```
In [3]: 1 df = pd.read_csv("7.1 heart-disease.csv")
        2 df.head()
```

```
Out[3]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
In [4]: 1 df.shape # (rows, columns)
```

```
Out[4]: (303, 14)
```

Data Exploration (exploratory data analysis or EDA)

The goal here is to find out more about the data and become a subject matter expert on the dataset you're working with.

1. What question(s) are you trying to solve?
2. what kind of data do we have and how do we treat different types?
3. What's missing from the data and how do you deal with it?
4. where are the outliers and why should you care about them?
5. how can you add, change or remove features to get more out of your data?

In [5]: 1 df.head()

Out[5]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [6]: 1 df.tail()

Out[6]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
299	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

In [7]: 1 df.target

Out[7]:

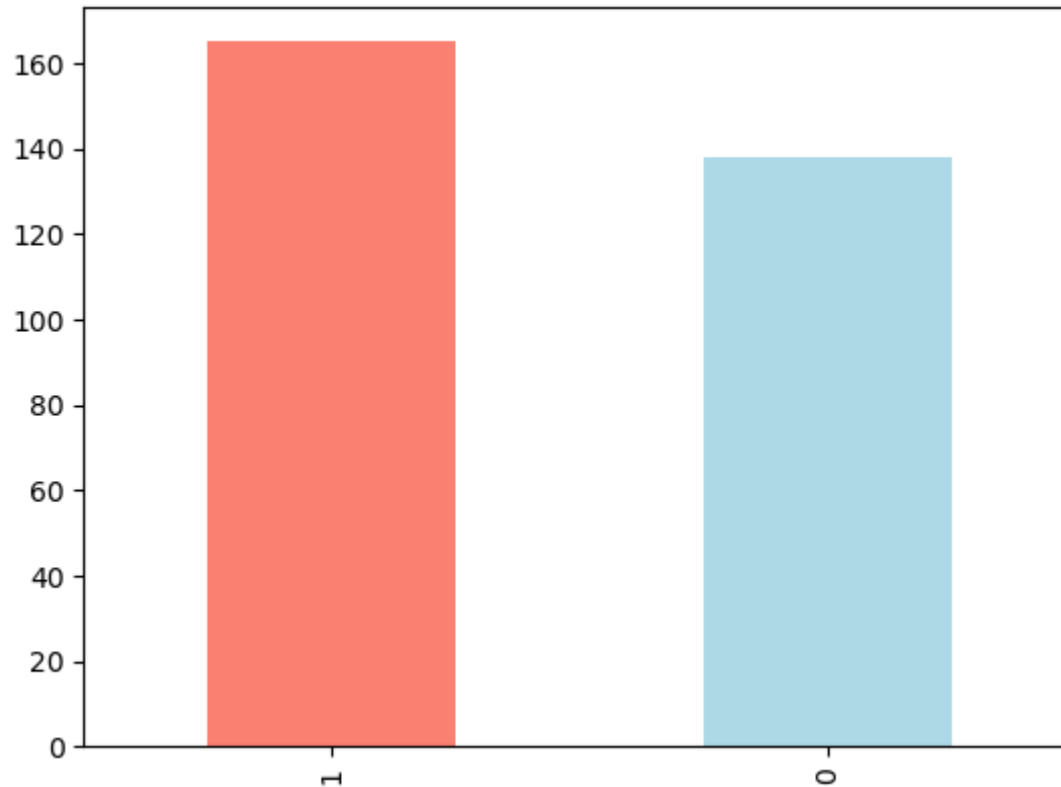
0	1
1	1
2	1
3	1
4	1
	..
298	0
299	0
300	0
301	0
302	0

Name: target, Length: 303, dtype: int64

```
In [8]: 1 # Let's find out how many of each class there  
2 df['target'].value_counts()
```

```
Out[8]: 1    165  
0     138  
Name: target, dtype: int64
```

```
In [9]: 1 df['target'].value_counts().plot(kind = 'bar', color = ['salmon', 'lightblue']);
```



In [10]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
---  -
 0   age         303 non-null   int64  
 1   sex         303 non-null   int64  
 2   cp          303 non-null   int64  
 3   trestbps    303 non-null   int64  
 4   chol        303 non-null   int64  
 5   fbs         303 non-null   int64  
 6   restecg     303 non-null   int64  
 7   thalach     303 non-null   int64  
 8   exang       303 non-null   int64  
 9   oldpeak     303 non-null   float64 
10  slope       303 non-null   int64  
11  ca          303 non-null   int64  
12  thal        303 non-null   int64  
13  target      303 non-null   int64  
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

In [11]:

```
1 # Are there any missing values?
2 df.isna().sum()
```

Out[11]:

```
age          0
sex          0
cp           0
trestbps     0
chol         0
fbs          0
restecg      0
thalach      0
exang        0
oldpeak      0
slope        0
ca           0
thal         0
target       0
dtype: int64
```

In [12]: 1 df.describe()

Out[12]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak
count	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
mean	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528053	149.646865	0.326733	1.039604
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000

Heart Disease frequency according to sex

In [13]: 1 df.sex.value_counts()

Out[13]:

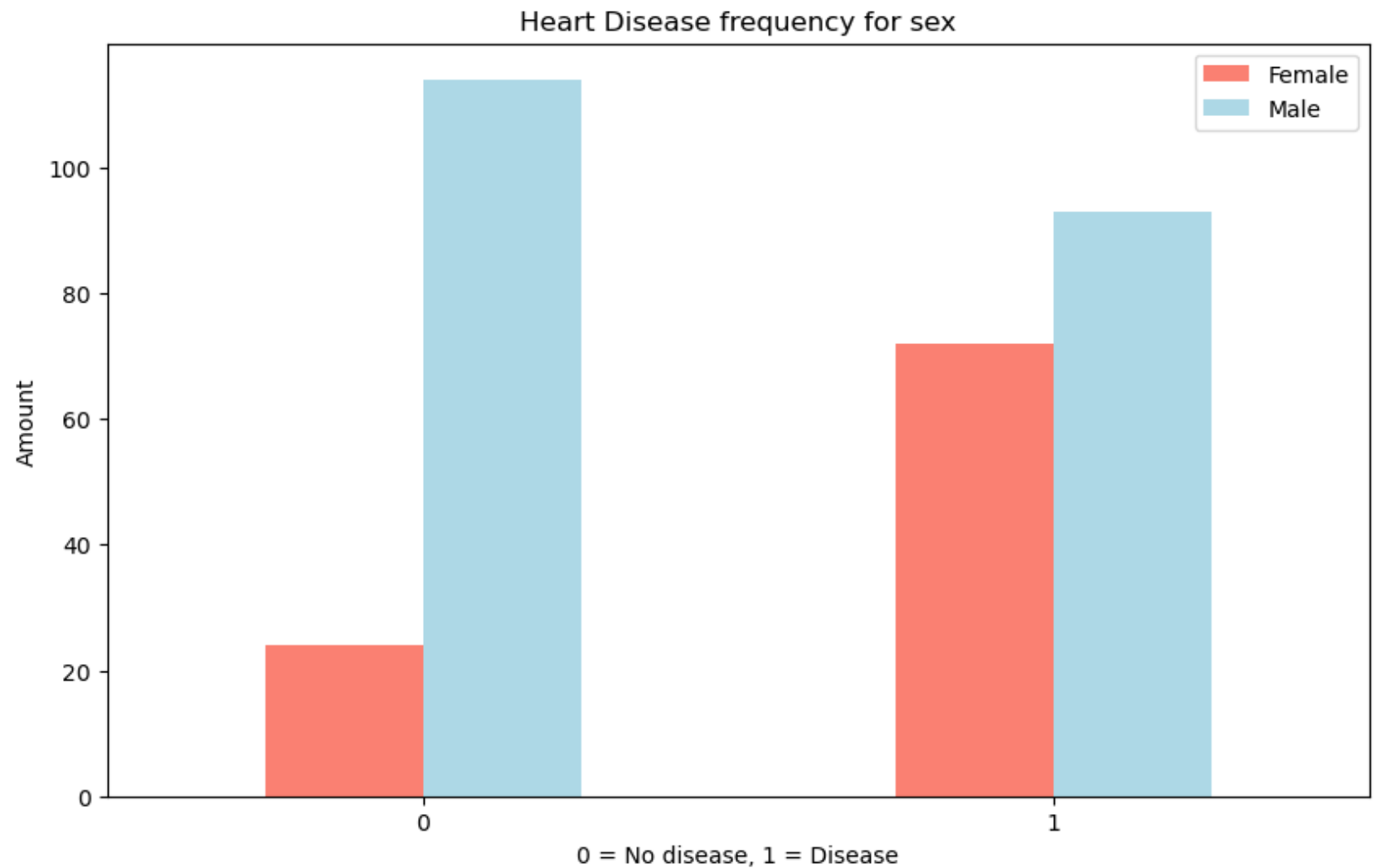
```
1    207
0     96
Name: sex, dtype: int64
```

In [14]: 1 # Compare target column with sex column
2 pd.crosstab(df.target, df.sex)

Out[14]:

	sex	0	1
target			
0	24	114	
1	72	93	


```
In [15]: 1 # Create a plot of crosstab
2 pd.crosstab(df.target, df.sex).plot(kind = 'bar',
3                                     figsize = (10, 6),
4                                     color = ['salmon', 'lightblue'])
5 plt.title('Heart Disease frequency for sex')
6 plt.xlabel('0 = No disease, 1 = Disease')
7 plt.ylabel('Amount')
8 plt.legend(['Female', 'Male'])
9 plt.xticks(rotation = 0);
```



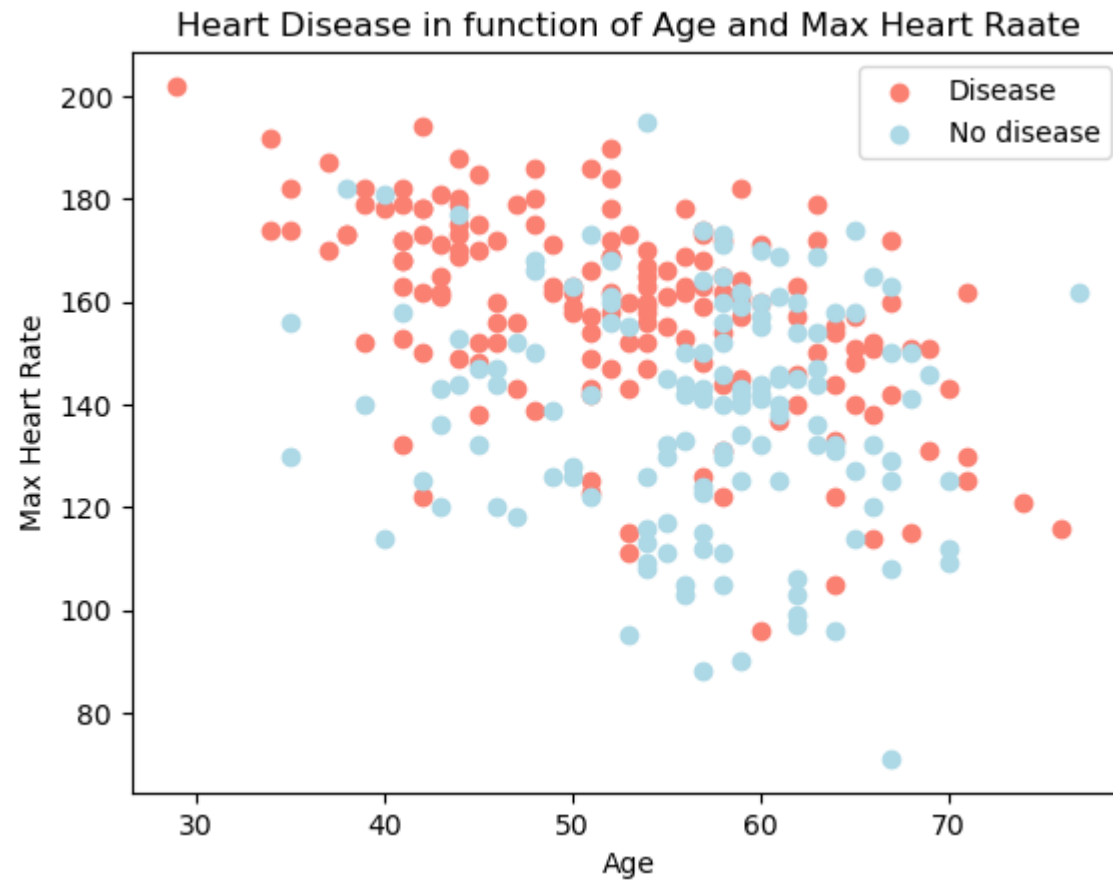
```
In [16]: 1 df['thalach'].value_counts()
```

```
Out[16]: 162    11
          160     9
          163     9
          152     8
          173     8
          ..
          202     1
          184     1
          121     1
          192     1
          90      1
          Name: thalach, Length: 91, dtype: int64
```

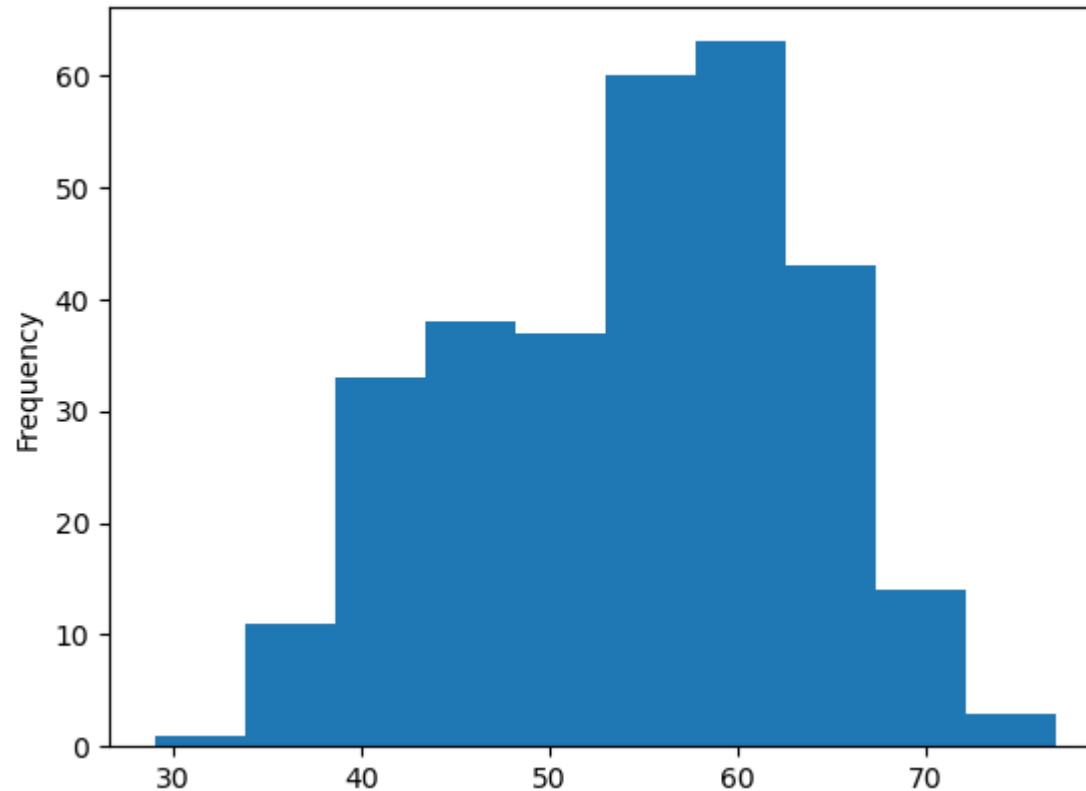
Age vs. Max Heart Rate for Heart Disease

In [17]:

```
1 # Create another figure
2 plt.figure = (10, 6)
3
4 # scatter with positive examples
5 plt.scatter(df.age[df.target == 1],
6             df.thalach[df.target == 1],
7             c = 'salmon', label = "Disease")
8
9 # scatter with negative examples
10 plt.scatter(df.age[df.target == 0],
11            df.thalach[df.target == 0],
12            c = 'lightblue', label = 'No disease');
13
14 # Add some helpful information
15 plt.title('Heart Disease in function of Age and Max Heart Raate')
16 plt.xlabel('Age')
17 plt.ylabel('Max Heart Rate')
18 plt.legend();
```



```
In [18]: 1 # check the distribution of the age column with a histogram  
2 df.age.plot.hist();
```



Heart Disease Frequency Chest pain Type

cp - chest pain type

- 0: Typical angina: chest pain related decrease blood supply to the heart
- 1: Atypical angina: chest pain not related to heart
- 2: Non-anginal pain: typically esophageal spasms (non heart related)
- 3: Asymptomatic: chest pain not showing signs of disease

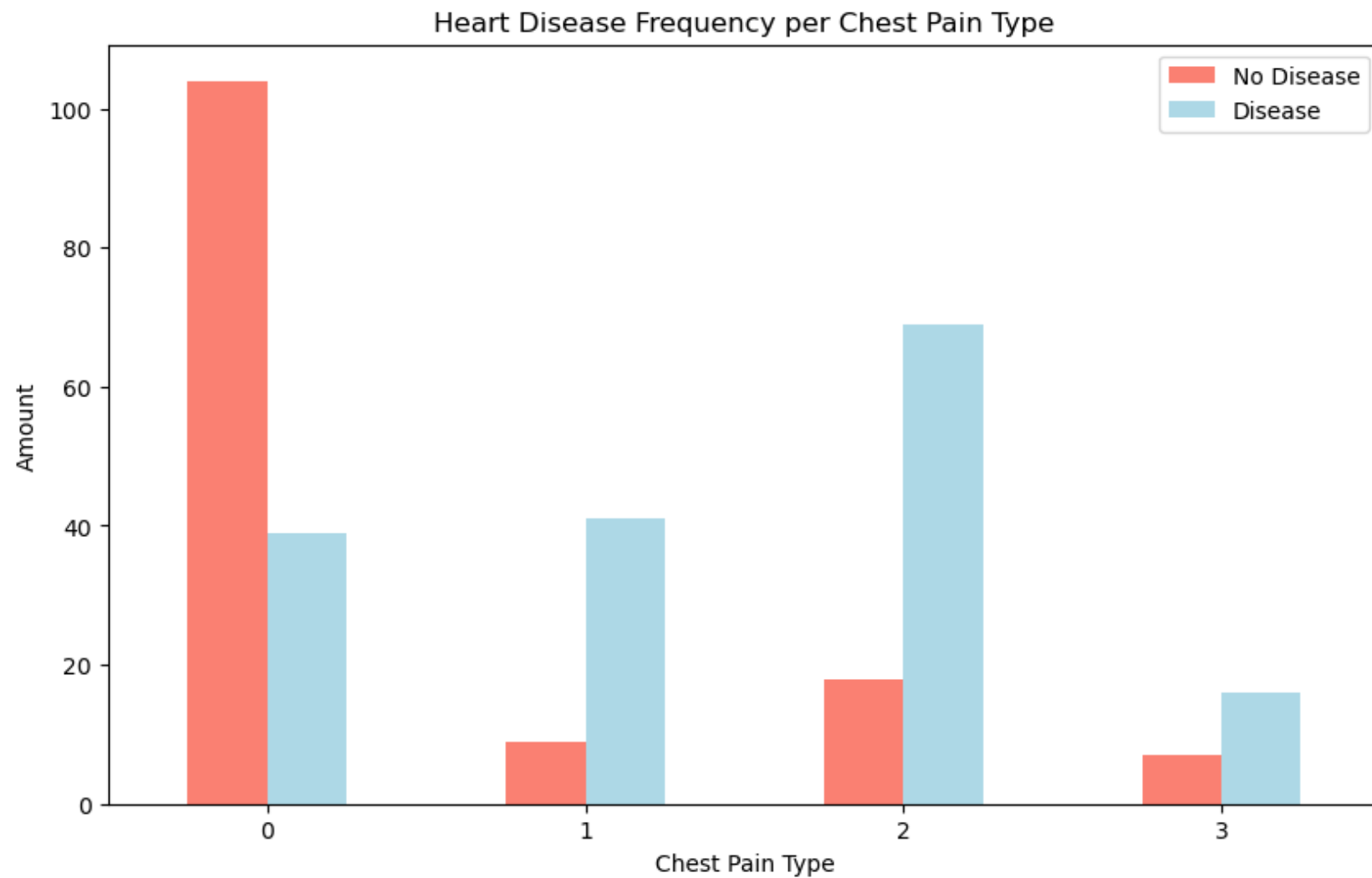
```
In [19]: 1 pd.crosstab(df.cp, df.target)
```

```
Out[19]:
```

target	0	1
cp		
0	104	39
1	9	41
2	18	69
3	7	16

In [20]:

```
1 # Make the crossbar more visuaal
2
3 pd.crosstab(df.cp, df.target).plot(kind = 'bar',
4                                     figsize = (10, 6),
5                                     color = ['salmon', 'lightblue'])
6
7 # Add soem communication
8 plt.title('Heart Disease Frequency per Chest Pain Type')
9 plt.xlabel('Chest Pain Type')
10 plt.ylabel('Amount')
11 plt.legend(['No Disease', 'Disease'])
12 plt.xticks(rotation = 0);
```



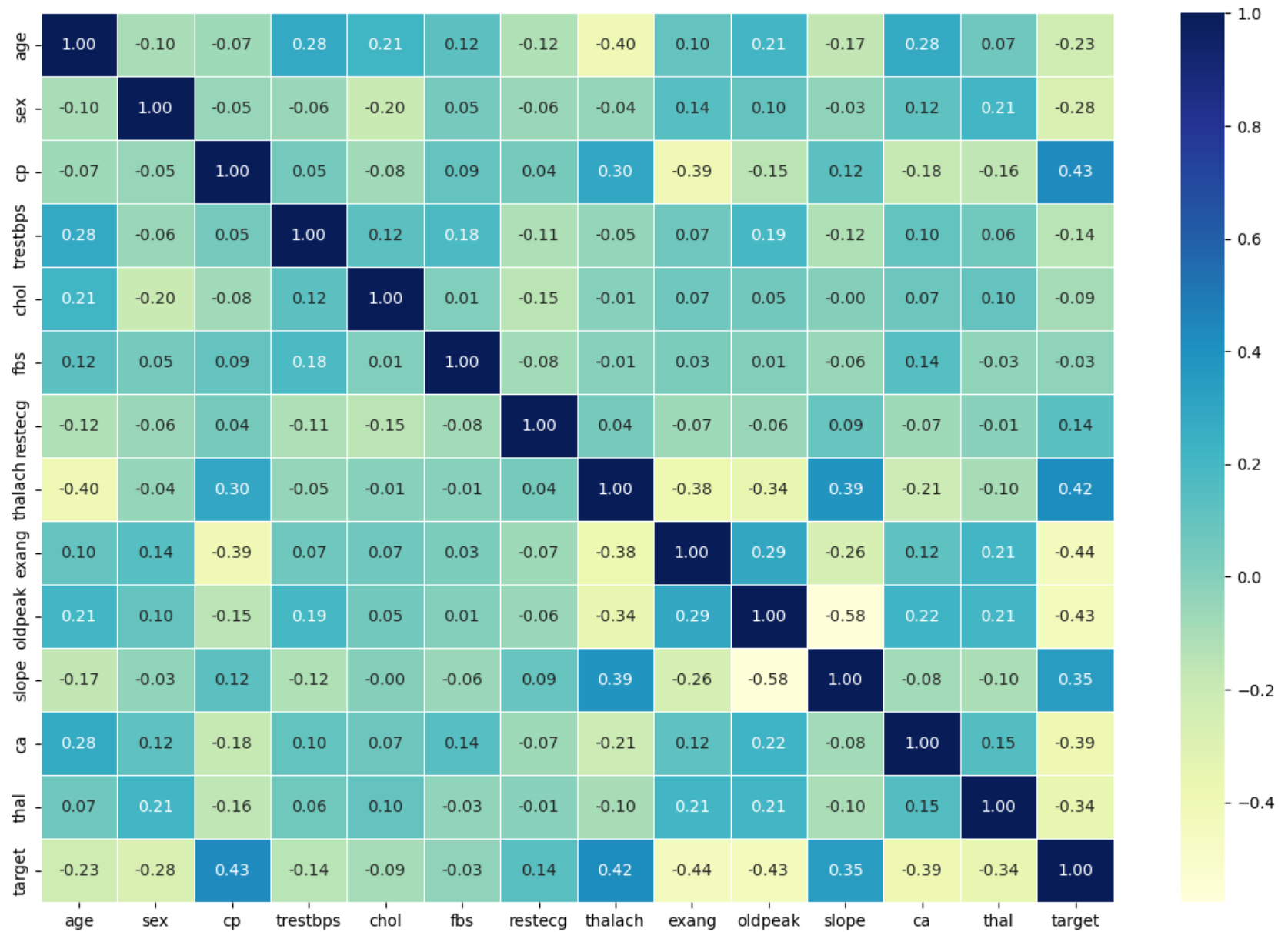
In [21]:

```
1 # Make a correlation matrix
2 df.corr()
```

Out[21]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	
age	1.000000	-0.098447	-0.068653	0.279351	0.213678	0.121308	-0.116211	-0.398522	0.096801	0.210013	-0.168814	0.27
sex	-0.098447	1.000000	-0.049353	-0.056769	-0.197912	0.045032	-0.058196	-0.044020	0.141664	0.096093	-0.030711	0.11
cp	-0.068653	-0.049353	1.000000	0.047608	-0.076904	0.094444	0.044421	0.295762	-0.394280	-0.149230	0.119717	-0.18
trestbps	0.279351	-0.056769	0.047608	1.000000	0.123174	0.177531	-0.114103	-0.046698	0.067616	0.193216	-0.121475	0.10
chol	0.213678	-0.197912	-0.076904	0.123174	1.000000	0.013294	-0.151040	-0.009940	0.067023	0.053952	-0.004038	0.07
fbs	0.121308	0.045032	0.094444	0.177531	0.013294	1.000000	-0.084189	-0.008567	0.025665	0.005747	-0.059894	0.13
restecg	-0.116211	-0.058196	0.044421	-0.114103	-0.151040	-0.084189	1.000000	0.044123	-0.070733	-0.058770	0.093045	-0.07
thalach	-0.398522	-0.044020	0.295762	-0.046698	-0.009940	-0.008567	0.044123	1.000000	-0.378812	-0.344187	0.386784	-0.21
exang	0.096801	0.141664	-0.394280	0.067616	0.067023	0.025665	-0.070733	-0.378812	1.000000	0.288223	-0.257748	0.11
oldpeak	0.210013	0.096093	-0.149230	0.193216	0.053952	0.005747	-0.058770	-0.344187	0.288223	1.000000	-0.577537	0.22
slope	-0.168814	-0.030711	0.119717	-0.121475	-0.004038	-0.059894	0.093045	0.386784	-0.257748	-0.577537	1.000000	-0.08
ca	0.276326	0.118261	-0.181053	0.101389	0.070511	0.137979	-0.072042	-0.213177	0.115739	0.222682	-0.080155	1.00
thal	0.068001	0.210041	-0.161736	0.062210	0.098803	-0.032019	-0.011981	-0.096439	0.206754	0.210244	-0.104764	0.15
target	-0.225439	-0.280937	0.433798	-0.144931	-0.085239	-0.028046	0.137230	0.421741	-0.436757	-0.430696	0.345877	-0.39

```
In [22]: 1 # Let's make our correlation matrix a little prettier
          2
          3 corr_matrix = df.corr()
          4 fig, ax = plt.subplots(figsize = (15, 10))
          5 ax = sns.heatmap(corr_matrix,
          6                   annot = True,
          7                   linewidths = 0.5,
          8                   fmt = '.2f',
          9                   cmap = 'YlGnBu');
```



5. Modelling

In [23]:

```
1 df.head()
```

Out[23]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [24]:

```
1 # split data into x and y
2
3 x = df.drop('target', axis = 1)
4
5 y = df['target']
```

In [25]:

```
1 x.head()
```

Out[25]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2

In [26]:

1	y
---	---

Out[26]:

0	1
1	1
2	1
3	1
4	1
	..
298	0
299	0
300	0
301	0
302	0

Name: target, Length: 303, dtype: int64

In [27]:

```
1 # split data into train and test sets
2 np.random.seed(42)
3
4 # split into train and test set
5 x_train, x_test, y_train, y_test = train_test_split(x,
6                                                     y,
7                                                     test_size=0.2)
```

In [28]: 1 x_train

Out[28]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
132	42	1	1	120	295	0	1	162	0	0.0	2	0	2
202	58	1	0	150	270	0	0	111	1	0.8	2	0	3
196	46	1	2	150	231	0	1	147	0	3.6	1	0	2
75	55	0	1	135	250	0	0	161	0	1.4	1	0	2
176	60	1	0	117	230	1	1	160	1	1.4	2	2	3
...
188	50	1	2	140	233	0	1	163	0	0.6	1	1	3
71	51	1	2	94	227	0	1	154	1	0.0	2	1	3
106	69	1	3	160	234	1	0	131	0	0.1	1	1	2
270	46	1	0	120	249	0	0	144	0	0.8	2	0	3
102	63	0	1	140	195	0	1	179	0	0.0	2	2	2

242 rows × 13 columns

In [29]: 1 y_train, len(y_train)

Out[29]:

```
(132    1
 202    0
 196    0
 75     1
 176    0
 ..
 188    0
 71     1
 106    1
 270    0
 102    1
 Name: target, Length: 242, dtype: int64,
 242)
```

Now we've got our data split into train and test sets, it's time to build a machine learning model.

We'll train it (find the patterns) on the training set.

And we'll test it (use the patterns) on the test set.

we're going to try 3 different machine learning models:

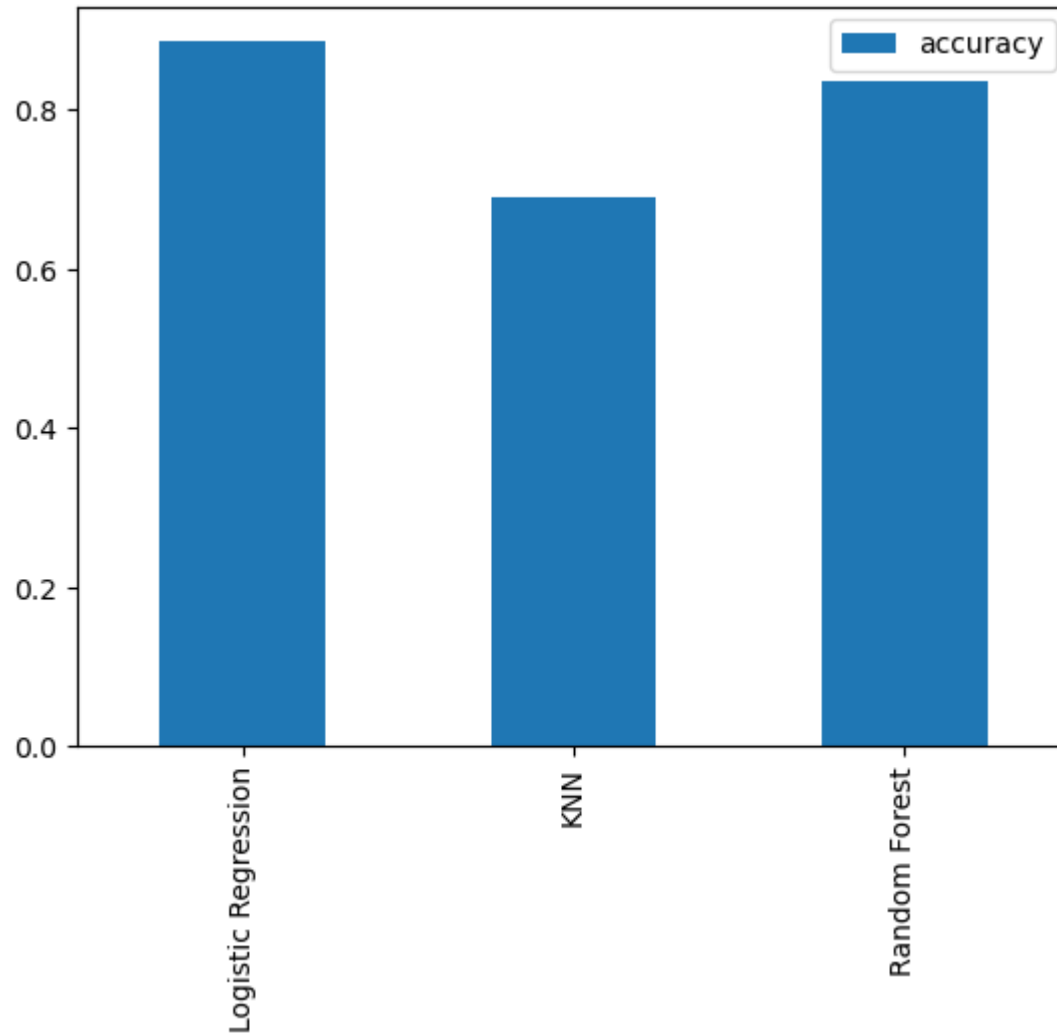
1. Logistic Regression
2. k-Nearest Neighbours Classifier
3. Random Forest Classifier

```
In [30]: 1 # Put models in a dictionary
2
3 models = {'Logistic Regression': LogisticRegression(),
4           'KNN': KNeighborsClassifier(),
5           'Random Forest': RandomForestClassifier()}
6
7 # Create a function to fit and score models
8 def fit_and_score(models, x_train, x_test, y_train, y_test):
9     '''
10     Fits and evaluate given machine learning models.
11     models: a dict of different scikit-learn machine learning models
12     x_train: training data (no labels)
13     x_test: testing data (no labels)
14     y_train: training labels
15     y_test: testing labels
16     '''
17     # set random seed
18     np.random.seed(42)
19     # make a dictionary to keep model scores
20     model_scores = {}
21     # Loop through models
22     for name, model in models.items():
23         # fit the model to the data
24         model.fit(x_train, y_train)
25         # Evaluate the model and append its score to model_score
26         model_scores[name] = model.score(x_test, y_test)
27     return model_scores
```

```
In [82]: 1 model_scores = fit_and_score(models = models,  
2                                     x_train = x_train,  
3                                     x_test = x_test,  
4                                     y_train = y_train,  
5                                     y_test = y_test)  
6 model_scores  
7 warnings.filterwarnings('ignore')
```

Model Comparison


```
In [32]: 1 model_compare = pd.DataFrame(model_scores, index=['accuracy'])  
2 model_compare.T.plot.bar();
```



Now we've got a baseline model.. and we know a model's first predictions aren't always what we should based our next steps off. What should do?

let's look at the following:

- Hyperparameter tuning

- features importance
- Confusion matrix
- cross-validation
- precision
- recall
- f1 score
- classification
- ROC Curve
- Area under the curve (AUC)

Hyperparameter tuning

```
In [64]: 1 # Let's tune KNN
          2
          3 train_scores = []
          4 test_scores = []
          5
          6 # create a list of different values for n_neighbors
          7 neighbors = range(1, 21)
          8
          9 # setup KNN instance
         10 knn = KNeighborsClassifier()
         11
         12 # Loop through different n_neighbors
         13 for i in neighbors:
         14     knn.set_params(n_neighbors=i)
         15
         16     # fit the algorithm
         17     knn.fit(x_train, y_train)
         18
         19     # update the training scores list
         20     train_scores.append(knn.score(x_train, y_train))
         21
         22     # update the test scores list
         23     test_scores.append(knn.score(x_test, y_test))
         24     import warnings
         25     warnings.filterwarnings('ignore')
```

```
In [34]: 1 train_scores
```

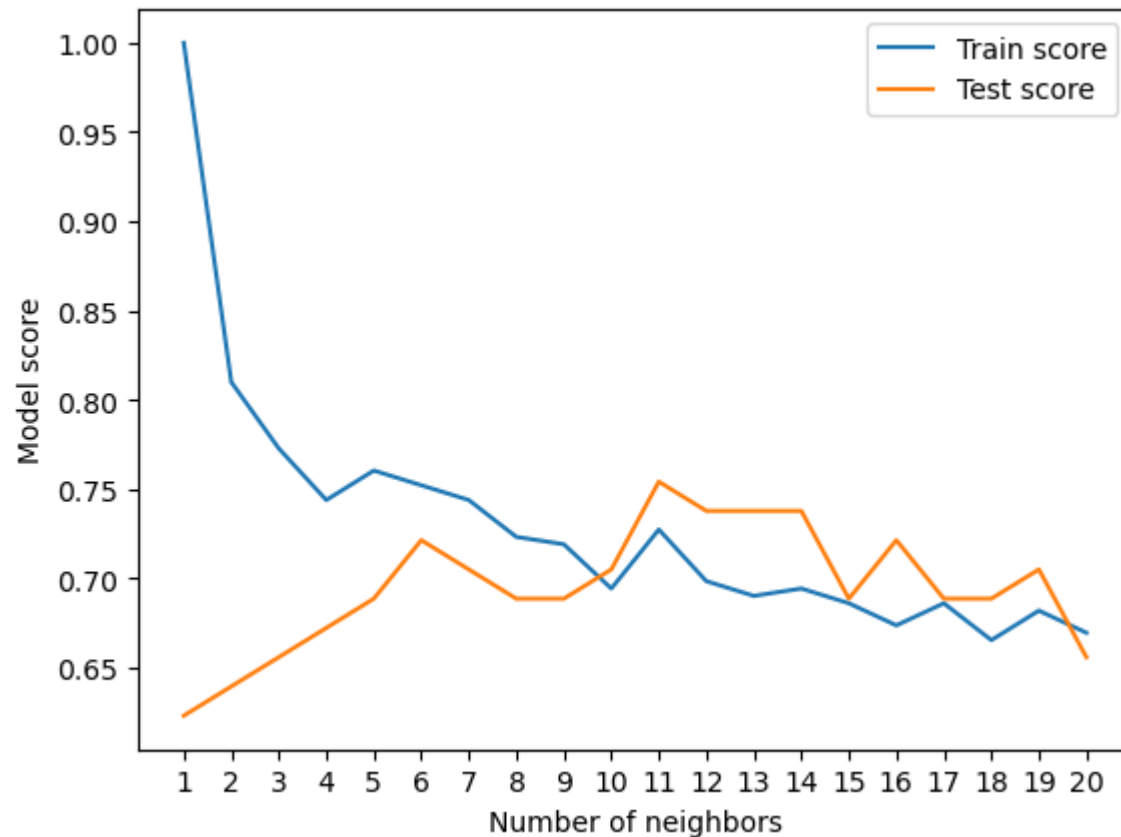
```
Out[34]: [1.0,  
0.8099173553719008,  
0.7727272727272727,  
0.743801652892562,  
0.7603305785123967,  
0.7520661157024794,  
0.743801652892562,  
0.7231404958677686,  
0.71900826446281,  
0.6942148760330579,  
0.7272727272727273,  
0.6983471074380165,  
0.6900826446280992,  
0.6942148760330579,  
0.6859504132231405,  
0.6735537190082644,  
0.6859504132231405,  
0.6652892561983471,  
0.6818181818181818,  
0.6694214876033058]
```

In [35]: 1 test_scores

Out[35]: [0.6229508196721312,
0.639344262295082,
0.6557377049180327,
0.6721311475409836,
0.6885245901639344,
0.7213114754098361,
0.7049180327868853,
0.6885245901639344,
0.6885245901639344,
0.7049180327868853,
0.7540983606557377,
0.7377049180327869,
0.7377049180327869,
0.7377049180327869,
0.6885245901639344,
0.7213114754098361,
0.6885245901639344,
0.6885245901639344,
0.7049180327868853,
0.6557377049180327]

```
In [36]: 1 plt.plot(neighbors, train_scores, label = 'Train score')
2 plt.plot(neighbors, test_scores, label = 'Test score')
3 plt.xticks(np.arange(1, 21, 1))
4 plt.xlabel('Number of neighbors')
5 plt.ylabel('Model score')
6 plt.legend()
7
8 print(f'Maximum KNN score on the test data: {max(test_scores) * 100:.2f}%')
```

Maximum KNN score on the test data: 75.41%



Hyperparameter tuning with RandomizedsearchCV

we're going to tune:

- LogisticRegression
- RandomForestClassifier

... using RandomizedSearchCV

```
In [37]: 1 # Create a hyperparameter grid for LogisticRegression
2 log_reg_grid = {'C': np.logspace(-4, 4, 20),
3               'solver': ['liblinear']}
4
5 # create a hyperparameter grid for RandomForestClassifier
6 rf_grid = {'n_estimators': np.arange(10, 1000, 50),
7           'max_depth': [None, 3, 5, 10],
8           'min_samples_split': np.arange(2, 20, 2),
9           'min_samples_leaf': np.arange(1, 20, 2)}
```

Now we've got hyperparameter grids setup for each of our models, RandomizedSearchCV...

```
In [38]: 1 # Tune LogisticRegression
2 np.random.seed(42)
3
4 # setup random hyperparameter search for LogisticRegression
5 rs_log_reg = RandomizedSearchCV(LogisticRegression(),
6                                param_distributions = log_reg_grid,
7                                cv=5,
8                                n_iter=20,
9                                verbose=True)
10
11 # Fit random hyperparameter search model for LogisticRegression
12 rs_log_reg.fit(x_train, y_train)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

```
Out[38]: RandomizedSearchCV(cv=5, estimator=LogisticRegression(), n_iter=20,
                          param_distributions={'C': array([1.00000000e-04, 2.63665090e-04, 6.95192796e-04, 1.83
298071e-03,
                          4.83293024e-03, 1.27427499e-02, 3.35981829e-02, 8.85866790e-02,
                          2.33572147e-01, 6.15848211e-01, 1.62377674e+00, 4.28133240e+00,
                          1.12883789e+01, 2.97635144e+01, 7.84759970e+01, 2.06913808e+02,
                          5.45559478e+02, 1.43844989e+03, 3.79269019e+03, 1.00000000e+04]),
                          'solver': ['liblinear']},
                          verbose=True)
```

In [39]: 1 rs_log_reg.best_params_

Out[39]: {'solver': 'liblinear', 'C': 0.23357214690901212}

In [40]: 1 rs_log_reg.score(x_test, y_test)

Out[40]: 0.8852459016393442

Now, we've tuned LogisticRegression(), let's do the same for RandomForestClassifier

```
In [41]: 1 # Setup random seed
2 np.random.seed(42)
3
4 # setup random hyperparameter search for RandomForestClassifier
5 rs_rf = RandomizedSearchCV(RandomForestClassifier(),
6                             param_distributions = rf_grid,
7                             cv=5,
8                             n_iter=20,
9                             verbose=True)
10
11 # fit random hyperparameter search model for RandomForestClassifier()
12 rs_rf.fit(x_train, y_train)
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

```
Out[41]: RandomizedSearchCV(cv=5, estimator=RandomForestClassifier(), n_iter=20,
                             param_distributions={'max_depth': [None, 3, 5, 10],
                             'min_samples_leaf': array([ 1,  3,  5,  7,  9, 11, 13, 15, 17, 1
9]),
                             'min_samples_split': array([ 2,  4,  6,  8, 10, 12, 14, 16, 1
8]),
                             'n_estimators': array([ 10,  60, 110, 160, 210, 260, 310, 360, 4
10, 460, 510, 560, 610,
                             660, 710, 760, 810, 860, 910, 960])},
                             verbose=True)
```

```
In [42]: 1 # find the best hyperparameters
         2 rs_rf.best_params_
```

```
Out[42]: {'n_estimators': 210,
          'min_samples_split': 4,
          'min_samples_leaf': 19,
          'max_depth': 3}
```

```
In [43]: 1 # Evaluation the randomized search RandomForestClassifier model
         2 rs_rf.score(x_test, y_test)
```

```
Out[43]: 0.8688524590163934
```

```
In [44]: 1 model_scores
```

```
Out[44]: {'Logistic Regression': 0.8852459016393442,
          'KNN': 0.6885245901639344,
          'Random Forest': 0.8360655737704918}
```

Hyperparameters using GridSearchCV

since our LogisticRegression model provides the best scores so far, we'll try and improve them again using GridSearchCV...


```
In [45]: 1 # Different hyperparamters for our LogisticRegression
2 log_reg_grid = {'C': np.logspace(-4, 4, 30),
3                 'solver': ['liblinear']}
4
5 # setup grid hyperparameter search for LogisticRegression
6 gs_log_reg = GridSearchCV(LogisticRegression(),
7                             param_grid = log_reg_grid,
8                             cv=5,
9                             verbose=True)
10
11 # fit grid hyperparameter search model
12 gs_log_reg.fit(x_train, y_train)
```

Fitting 5 folds for each of 30 candidates, totalling 150 fits

```
Out[45]: GridSearchCV(cv=5, estimator=LogisticRegression(),
                      param_grid={'C': array([1.00000000e-04, 1.88739182e-04, 3.56224789e-04, 6.72335754e-04,
1.26896100e-03, 2.39502662e-03, 4.52035366e-03, 8.53167852e-03,
1.61026203e-02, 3.03919538e-02, 5.73615251e-02, 1.08263673e-01,
2.04335972e-01, 3.85662042e-01, 7.27895384e-01, 1.37382380e+00,
2.59294380e+00, 4.89390092e+00, 9.23670857e+00, 1.74332882e+01,
3.29034456e+01, 6.21016942e+01, 1.17210230e+02, 2.21221629e+02,
4.17531894e+02, 7.88046282e+02, 1.48735211e+03, 2.80721620e+03,
5.29831691e+03, 1.00000000e+04]),
                      'solver': ['liblinear']},
                      verbose=True)
```

```
In [46]: 1 # check the best hyperparameters
2 gs_log_reg.best_params_
```

```
Out[46]: {'C': 0.20433597178569418, 'solver': 'liblinear'}
```

```
In [47]: 1 # Evaluate the grid search Logisticregression model
2 gs_log_reg.score(x_test, y_test)
```

```
Out[47]: 0.8852459016393442
```

Evaluating our tuned machine learning classifier, beyond accuracy

- ROC curve and AUC score
- Confusion matrix
- Classification report
- Precision
- Recall
- F1-score

...and it would be great if cross-validation was used where possible.

To make comparison and evaluate our model first we need to make predictions

```
In [48]: 1 # make predictions with tuned model
          2 y_preds = gs_log_reg.predict(x_test)
```

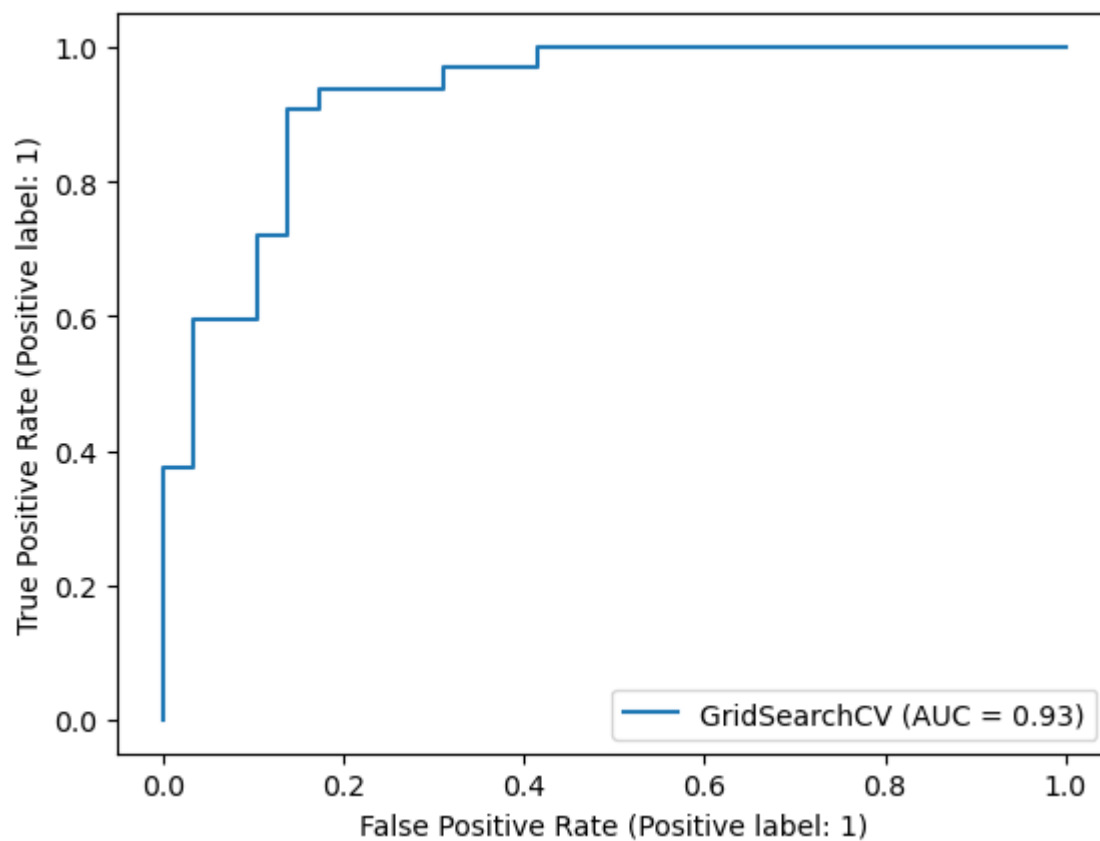
```
In [49]: 1 y_preds
```

```
Out[49]: array([0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0,
                0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
                1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0], dtype=int64)
```

```
In [50]: 1 y_test
```

```
Out[50]: 179    0
          228    0
          111    1
          246    0
           60    1
           ..
          249    0
          104    1
          300    0
          193    0
          184    0
          Name: target, Length: 61, dtype: int64
```

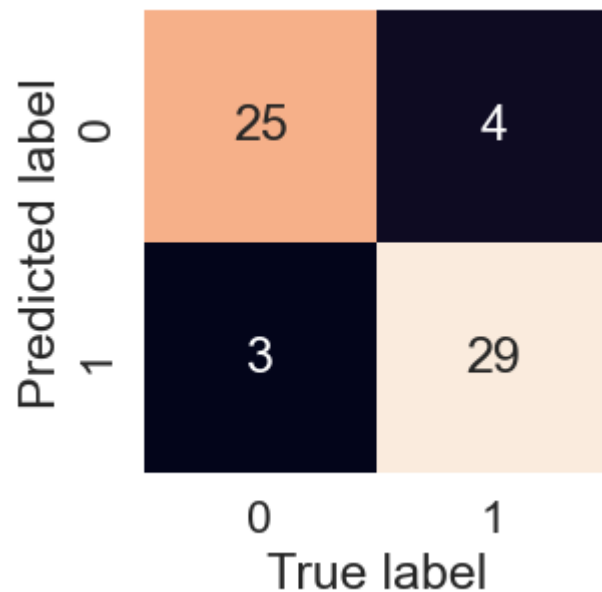
```
In [51]: 1 # plot ROC curve and calculate and calculate AUC metrics  
2 plot_roc_curve(gs_log_reg, x_test, y_test);
```



```
In [52]: 1 # confusion matrix  
2 print(confusion_matrix(y_test, y_preds))
```

```
[[25  4]  
 [ 3 29]]
```

```
In [53]: 1 sns.set(font_scale = 1.5)
2
3 def plot_conf_mat(y_test, y_preds):
4     '''
5     Plots a nice looking confusion matrix using seaborn's heatmap()
6     '''
7     fig, ax = plt.subplots(figsize=(3, 3))
8     ax = sns.heatmap(confusion_matrix(y_test, y_preds),
9                     annot=True,
10                    cbar=False)
11     plt.xlabel('True label')
12     plt.ylabel('Predicted label')
13 plot_conf_mat(y_test, y_preds)
```



Now we've got a ROC curve, an AUC metric and a confusion matrix, let's get classification report as well as cross-validation precision, recall and f1-score.

In [54]: 1 `print(classification_report(y_test, y_preds))`

	precision	recall	f1-score	support
0	0.89	0.86	0.88	29
1	0.88	0.91	0.89	32
accuracy			0.89	61
macro avg	0.89	0.88	0.88	61
weighted avg	0.89	0.89	0.89	61

calculate evaluation metrics using cross-validation

we're going to calculate accuracy, precision, recall and f1-score of our model using cross-validation and to do so we'll be using 'cross_val_score().'

In [55]: 1 `# check best hyperparameters`
2 `gs_log_reg.best_params_`

Out[55]: {'C': 0.20433597178569418, 'solver': 'liblinear'}

In [56]: 1 `# create a new classifier with best parameters`
2
3 `clf = LogisticRegression(C=0.20433597178569418,`
4 `solver= 'liblinear')`

In [57]: 1 `# cross-validated accuracy`
2
3 `cv_acc = cross_val_score(clf,`
4 `x,`
5 `y,`
6 `cv=5,`
7 `scoring='accuracy')`
8 `cv_acc1 = (np.mean(cv_acc))`
9 `cv_acc1`

Out[57]: 0.8446994535519124

```
In [58]: 1 # cross-validated precision
2 cv_precision = cross_val_score(clf,
3                                 x,
4                                 y,
5                                 cv=5,
6                                 scoring='precision')
7 cv_precision1 = np.mean(cv_precision)
8 cv_precision1
```

Out[58]: 0.8207936507936507

```
In [59]: 1 # cross-validated recall
2 cv_recall = cross_val_score(clf,
3                              x,
4                              y,
5                              cv=5,
6                              scoring='recall')
7 cv_recall1 = np.mean(cv_recall)
8 cv_recall1
```

Out[59]: 0.9212121212121213

```
In [60]: 1 # cross-validated f1-score
2
3 cv_f1_score = cross_val_score(clf,
4                                x,
5                                y,
6                                cv=5,
7                                scoring='f1')
8 cv_f1_score1 = np.mean(cv_recall)
9 cv_f1_score1
```

Out[60]: 0.9212121212121213

```
In [61]: 1 # visualize cross-validated-metrics
2 cv_metrics = pd.DataFrame({'Accuracy': cv_acc1,
3                             'Precision': cv_precision1,
4                             'Recall': cv_recall1,
5                             'F1': cv_f1_score1},
6                             index=[0] )
7 cv_metrics
```

```
Out[61]:
```

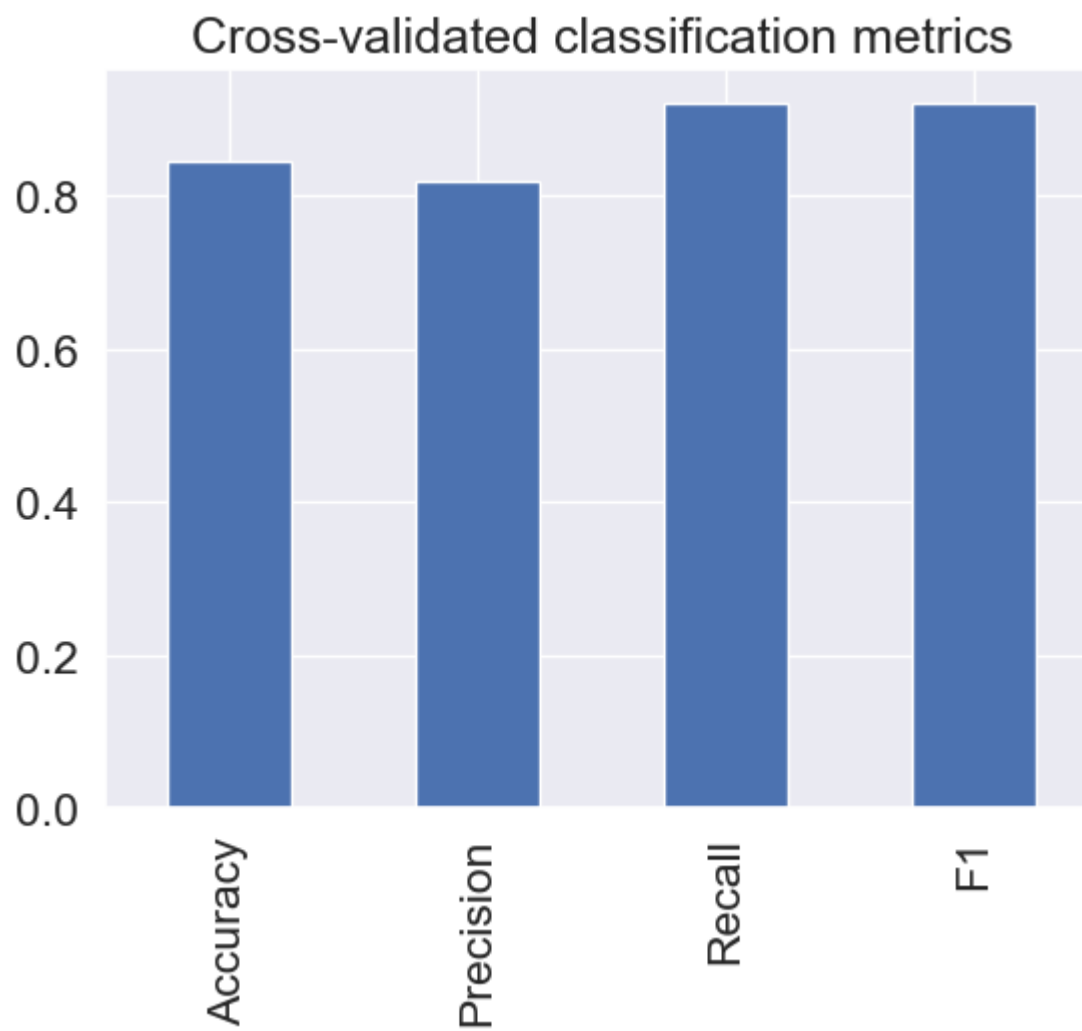
	Accuracy	Precision	Recall	F1
0	0.844699	0.820794	0.921212	0.921212

```
In [62]: 1 # visualize cross-validated-metrics
2 cv_metrics = pd.DataFrame({'Accuracy': cv_acc,
3                             'Precision': cv_precision,
4                             'Recall': cv_recall,
5                             'F1': cv_f1_score},
6                             )
7 cv_metrics
```

```
Out[62]:
```

	Accuracy	Precision	Recall	F1
0	0.819672	0.775000	0.939394	0.849315
1	0.901639	0.885714	0.939394	0.911765
2	0.868852	0.857143	0.909091	0.882353
3	0.883333	0.861111	0.939394	0.898551
4	0.750000	0.725000	0.878788	0.794521

```
In [63]: 1 # visualize cross-validated-metrics
2 cv_metrics = pd.DataFrame({'Accuracy': cv_acc1,
3                             'Precision': cv_precision1,
4                             'Recall': cv_recall1,
5                             'F1': cv_f1_score1},
6                             index=[0] )
7 cv_metrics.T.plot.bar(title='Cross-validated classification metrics',
8                        legend = False);
```



Feature importance

feature importance is another way of asking, 'which features contributed most to the outcome of the model and how did they contribute?'

Finding feature importance is different for each machine learning models. One way to find feature importance is to search for '(Model NAME)' feature importance.

let's find the feature importance for our LogisticRegression model...which gave the best result out of the three models

In [65]:

```
1 df.head()
```

Out[65]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

In [66]:

```
1 # Fit an instance of LogisticRegression
2 gs_log_reg.best_params_
3
4
```

Out[66]: {'C': 0.20433597178569418, 'solver': 'liblinear'}

In [67]:

```
1 clf = LogisticRegression(C=0.20433597178569418,
2                             solver='liblinear')
3 clf.fit(x_train, y_train);
```

```
In [68]: 1 # check coef_
        2 clf.coef_
```

```
Out[68]: array([[ 0.00316728, -0.86044655,  0.66067042, -0.01156993, -0.00166374,
  0.04386109,  0.31275847,  0.02459361, -0.60413083, -0.56862804,
  0.4505163 , -0.63609898, -0.67663378]])
```

```
In [74]: 1 df.head()
```

```
Out[74]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1

```
In [73]: 1 # match coef's of features to columns
        2 feature_dict = dict(zip(df.columns, list(clf.coef_[0])))
        3 feature_dict
```

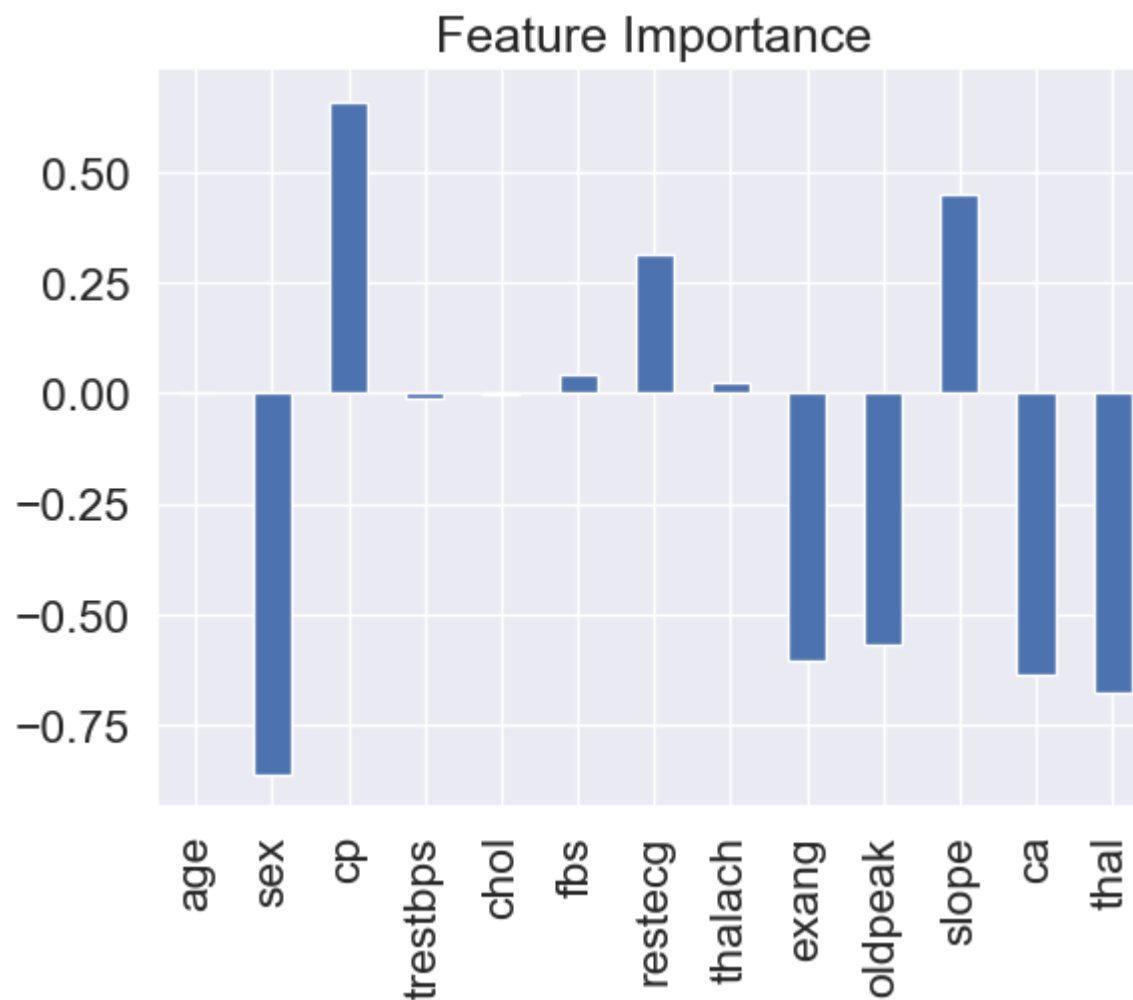
```
Out[73]: {'age': 0.0031672809701328098,
'sex': -0.8604465542018854,
'cp': 0.6606704161071124,
'trestbps': -0.011569931839584581,
'chol': -0.0016637442846940298,
'fbs': 0.043861090099753856,
'restecg': 0.3127584688139112,
'thalach': 0.024593614036076978,
'exang': -0.6041308274033194,
'oldpeak': -0.5686280446250761,
'slope': 0.45051629703183155,
'ca': -0.6360989766185763,
'thal': -0.6766337834775279}
```

```
In [75]: 1 # visualize featurte importance
2
3 feature_df = pd.DataFrame(feature_dict, index=[0])
4 feature_df
```

```
Out[75]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	t
0	0.003167	-0.860447	0.66067	-0.01157	-0.001664	0.043861	0.312758	0.024594	-0.604131	-0.568628	0.450516	-0.636099	-0.6766

```
In [77]: 1 feature_df.T.plot.bar(title = 'Feature Importance', legend = False);
```



```
In [78]: 1 pd.crosstab(df['sex'], df['target'])
```

```
Out[78]:
```

	target	0	1
sex			
0	24	72	
1	114	93	

```
In [79]: 1 pd.crosstab(df['slope'], df['target'])
```

```
Out[79]:
```

	target	0	1
slope			
0	12	9	
1	91	49	
2	35	107	

slope - the slope of the peak exercise ST segment

- 0: Upsloping: better heart rate with exercise (uncommon)
- 1: Flatsloping: minimal change (typical healthy heart)
- 2: Downsloping: signs of unhealthy heart

6. Experimentation

if you haven't hit your evaluation metric yet... ask yourself...

- could you collect more data?
- could you try a better model? like CatBoost or XGBoost?
- could you improve the current models? (beyond we've done so far)
- if your model is good enough (you have hit your evaluation metrics) how would you export it and share it with others?

In []:

1